

```

1. for j = 1 : m
2.   w_j = Aq_j
3.   for i = 1 : j
4.     h_ij = (w_j, q_i)
5.     w_j = w_j - h_ij q_i           % subtract projection
6.   endfor
7.   h_{j+1,j} = ||w_j||_2
8.   if h_{j+1,j} == 0
9.     stop
10.  endif
11.  q_{j+1} = w_j / h_{j+1,j}       % normalize
12. endfor

```

□

Lemma 5.4. Computation of an orthonormal basis by Arnoldi's method. *If $\dim K_m(\underline{q}_1, A) = m$, then Arnoldi's method computes an orthonormal basis $\{\underline{q}_1, \dots, \underline{q}_m\}$ of $K_m(\underline{q}_1, A)$.*

Proof. The vectors $\underline{q}_1, \dots, \underline{q}_m$ are orthonormal by construction: orthogonal by line 3 – 6 and normalized by line 7 and 11. One has to show that they belong all to $K_m(\underline{q}_1, A)$. This statement will follow from the fact that each vector \underline{q}_j is of the form $p_{j-1}(A)\underline{q}_1$, where p_{j-1} is a polynomial of degree $j-1$. The proof is done by induction. For $j=1$, one has $\underline{q}_1 = p_0(A)\underline{q}_1$ such that $p_0(t) = 1$. Assume, the statement is true for $j \leq k$. One has by using first line 11, which implies $h_{k+1,k} \neq 0$, then lines 2 – 6, and finally the assumption of the induction

$$\begin{aligned}
h_{k+1,k}\underline{q}_{k+1} &= \underline{w}_k = A\underline{q}_k - \sum_{i=1}^k h_{ik}\underline{q}_i \\
&= Ap_{k-1}(A)\underline{q}_1 - \sum_{i=1}^k h_{ik}p_{i-1}(A)\underline{q}_1 = p_k(A)\underline{q}_1.
\end{aligned} \tag{5.1}$$

Division by $h_{k+1,k}$ finishes the proof. ■

Remark 5.5. Factorization of the system matrix. Denote

$$\begin{aligned}
Q_m &= (\underline{q}_1, \underline{q}_2, \dots, \underline{q}_m) \in \mathbb{R}^{n \times m}, \\
H_m &= \begin{pmatrix} h_{11} & h_{12} & \cdots & h_{1,m-1} & h_{1m} \\ h_{21} & h_{22} & \cdots & h_{2,m-1} & h_{2m} \\ 0 & h_{32} & \cdots & h_{3,m-1} & h_{3m} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & h_{m,m-1} & h_{mm} \\ 0 & 0 & \cdots & 0 & h_{m+1,m} \end{pmatrix} \in \mathbb{R}^{(m+1) \times m}.
\end{aligned}$$

A matrix of this form, i.e., $h_{ij} = 0$ for $i > j+1$, is called (upper) Hessenberg⁶ matrix. It follows from (5.1) that

$$A\underline{q}_k = \underline{q}_{k+1}h_{k+1,k} + \sum_{i=1}^k \underline{q}_i h_{ik} = \sum_{i=1}^{k+1} \underline{q}_i h_{ik}, \quad k = 1, \dots, m, \tag{5.2}$$

such that one gets for Arnoldi's method the compact representation

$$AQ_m = Q_{m+1}H_m. \tag{5.3}$$

□

⁶ Karl Hessenberg (1904 – 1959)

Remark 5.6. Initial vector in Krylov subspace methods. In the Krylov subspace methods, $\underline{r}^{(0)} / \|\underline{r}^{(0)}\|_2$ plays the role of \underline{q}_1 in Arnoldi's method. \square

Remark 5.7. Principle approach for minimizing the residual. The goal of the methods presented in this section is to minimize the Euclidean norm of the residual. One has

$$\begin{aligned} \|\underline{r}^{(k)}\|_2 &= \|\underline{b} - A\underline{x}^{(k)}\|_2 = \|\underline{r}^{(0)} + A\underline{x}^{(0)} - A\underline{x}^{(k)}\|_2 \\ &= \|\underline{r}^{(0)} - A(\underline{x}^{(k)} - \underline{x}^{(0)})\|_2 \end{aligned}$$

with $\underline{x}^{(k)} - \underline{x}^{(0)} \in K_k(\underline{r}^{(0)}, A)$, see Remark 4.9. Since the vectors $\{\underline{q}_1, \dots, \underline{q}_k\}$ computed with Arnoldi's method form a generating system of $K_k(\underline{r}^{(0)}, A)$, it is

$$\underline{x}^{(k)} - \underline{x}^{(0)} = \sum_{i=1}^k z_i \underline{q}_i = Q_k \underline{z} \quad (5.4)$$

with $\underline{z} = (z_1, \dots, z_k)^T$, $Q_k = (\underline{q}_1, \dots, \underline{q}_k)$. Using (5.4), (5.3), $Q_{k+1} \underline{e}_1 = \underline{q}_1 = \underline{r}^{(0)} / \|\underline{r}^{(0)}\|_2$, and the fact that the Euclidean norm is invariant under orthonormal transformations, one obtains

$$\begin{aligned} \|\underline{r}^{(k)}\|_2^2 &= \|\underline{r}^{(0)} - AQ_k \underline{z}\|_2^2 = \|\underline{r}^{(0)} - Q_{k+1} H_k \underline{z}\|_2^2 \\ &= \left\| \|\underline{r}^{(0)}\|_2 Q_{k+1} \underline{e}_1 - Q_{k+1} H_k \underline{z} \right\|_2^2 = \left\| Q_{k+1} \left(\|\underline{r}^{(0)}\|_2 \underline{e}_1 - H_k \underline{z} \right) \right\|_2^2 \\ &= \left\| \|\underline{r}^{(0)}\|_2 \underline{e}_1 - H_k \underline{z} \right\|_2^2. \end{aligned}$$

The minimizer of the residual is obtained by solving the least squares problem

$$\min_{\underline{z} \in \mathbb{R}^k} \left\| \|\underline{r}^{(0)}\|_2 \underline{e}_1 - H_k \underline{z} \right\|_2^2. \quad (5.5)$$

This problem possesses k unknowns z_1, \dots, z_k and the vector that has to be minimized has $k+1$ components. It can be solved, e.g., with a QR factorization of H_k , see lecture notes Numerical Mathematics I. Let $\underline{z}^{(k)}$ be a solution of this problem, then the next iterate of the Krylov subspace method is, compare (5.4)

$$\underline{x}^{(k)} = \underline{x}^{(0)} + Q_k \underline{z}^{(k)}. \quad (5.6)$$

This algorithm is called GMRES (generalized minimal residual). It has been proposed the first time in Saad & Schultz (1986). \square

Theorem 5.8. Properties of GMRES.

- i) In the case that Arnoldi's method has an early breakdown, i.e., $h_{l+1,l} = 0$, then $\dim K_k(\underline{r}^{(0)}, A) = l < k$ and $\underline{r}^{(l)} = \underline{0}$. Hence $A\underline{x}^{(l)} = \underline{b}$.*
- ii) The iterate $\underline{x}^{(k)} = \underline{x}^{(0)} + Q_k \underline{z}^{(k)}$ is uniquely determined.*
- iii) It holds*

$$\|\underline{r}^{(k)}\|_2 \leq \|\underline{r}^{(k-1)}\|_2, \quad k = 1, 2, 3, \dots$$

Proof. *i)* The breakdown of Arnoldi's method after l steps, $h_{l+1,l} = 0$, is equivalent to $\underline{w}_l = \underline{0}$, see line 8. It follows from (5.1) that

$$A\underline{q}_l = \sum_{i=1}^l h_{il} \underline{q}_i,$$

where $\underline{q}_1 = \underline{r}^{(0)} / \|\underline{r}^{(0)}\|_2$ in the case of GMRES. Since, in contrast to (5.2) (upper limit of the summation), $A\underline{q}_l$ is a linear combination of the already computed l basis vectors, one has $\dim K_{l+1}(\underline{r}^{(0)}, A) = \dim K_l(\underline{r}^{(0)}, A)$. One obtains by induction

$$\dim K_k(\underline{r}^{(0)}, A) = \dim K_l(\underline{r}^{(0)}, A) \text{ for } k \geq l.$$

Using matrix notations, Arnoldi's method gives in this case

$$AQ_l = Q_l \tilde{H}_l \text{ with } \tilde{H}_l = \begin{pmatrix} h_{11} & h_{12} & \cdots & h_{1l} \\ h_{21} & h_{22} & \cdots & h_{2l} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & h_{l,l-1} & h_{ll} \end{pmatrix} \in \mathbb{R}^{l \times l}, \quad Q_l \in \mathbb{R}^{n \times l}.$$

Since A is non-singular and $\text{rank}(Q_l) = l$, one has $\text{rank}(AQ_l) = l$. Consequently, it is $\text{rank}(Q_l \tilde{H}_l) = l$ and $\text{rank}(\tilde{H}_l) = l$ and \tilde{H}_l is invertible. In the same way as in Remark 5.7, one obtains

$$\|\underline{r}^{(l)}\|_2^2 = \min_{\underline{z} \in \mathbb{R}^l} \left\| \|\underline{r}^{(0)}\|_2 \underline{e}_1 - \tilde{H}_l \underline{z} \right\|_2^2.$$

The minimizer is given by $\underline{z}^{(l)} = \tilde{H}_l^{-1} \left(\|\underline{r}^{(0)}\|_2 \underline{e}_1 \right)$ which gives $\|\underline{r}^{(l)}\|_2 = 0$.

ii) If $\text{rank}(H_k) = k$, the minimizer of (5.5) is unique (theory of least squares problems, see Numerical Mathematics I). If $\text{rank}(H_k) < k$, then $\underline{x}^{(k)} = \underline{x}$, see i).

iii) The set in which the minimizer is computed becomes larger since the inclusion $K_k(\underline{r}^{(0)}, A) \supseteq K_{k-1}(\underline{r}^{(0)}, A)$ holds. ■

Remark 5.9. Implementational issues.

- The GMRES process consists in principle of two steps:

1. computing the orthonormal basis of $K_k(\underline{r}^{(0)}, A)$,

2. solving the least squares problem (5.5) to find the minimizer of the residual with a standard method.

In the practical use of GMRES, Step 2 is performed only at the end of the iteration. Thus, the iterate $\underline{x}^{(k)}$ is not directly available. It is computed in a post-processing step. However, there is an elegant and inexpensive way to compute $\|\underline{r}^{(k)}\|_2$ without having access to $\underline{x}^{(k)}$, see Saad (2003). With $\|\underline{r}^{(k)}\|_2$, one can control the iterative process.

For concrete ways to implement GMRES, it is referred to Saad & Schultz (1986); Saad (2003).

- Each step of GMRES requires one matrix-vector multiplication, line 2 of Arnoldi's method.
- In exact arithmetic, GMRES terminates with the solution in at most n steps. This property is, however, of no practical use for large n .
- From the practical point of view, the greatest problem of GMRES is that one has to store the basis $\{\underline{q}_1, \dots, \underline{q}_k\}$ of $K_k(\underline{r}^{(0)}, A)$, see lines 3 – 6 of Arnoldi's method. Thus, with every new iteration, one has to store an additional vector. This situation is called long recurrence.

In practice, one prescribes a maximal dimension m of the Krylov subspace. After m iterations, GMRES is stopped with the iterate $\underline{x}^{(m)}$. If $\underline{x}^{(m)}$ is not yet sufficiently close to the solution, GMRES is started from the beginning with $\underline{x}^{(0)} = \underline{x}^{(m)}$. This approach is called GMRES(m) (with restart). An optimal choice of m is in general an unresolved problem. Often $m \in [5, 50]$ is used.

GMRES(m) might also fail to converge, see Saad & Schultz (1986) for the simple example

$$A = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}, \quad \underline{b} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \underline{x}^{(0)} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

GMRES converges in two steps whereas GMRES(1) computes the stationary sequence $\underline{x}^{(1)} = \underline{x}^{(0)}$, $\underline{x}^{(2)} = \underline{x}^{(1)} = \underline{x}^{(0)}$ and so on. Despite of the possibility to fail, GMRES(m) is one of the most popular and best performing iterative methods for solving linear systems of equations with non-symmetric matrix. □

5.2 Symmetric Matrices

Remark 5.10. Goal. Arnoldi's method and the minimization of the residual in $K_k(r^{(0)}, A)$ will be studied in the special case that A is symmetric. The most important result in this case will be that it is not necessary to store the basis of $K_k(r^{(0)}, A)$. It suffices to store a fixed number of only few basis vectors. Thus, the memory requirements do not increase in the course of the iteration and the most important problem of using GMRES vanishes. \square

Remark 5.11. Arnoldi's method revisited. First, Arnoldi's method is revisited. From the general relation (5.3), it follows by the orthonormality of the columns of Q_k and Q_{k+1} that

$$Q_k^T A Q_k = Q_k^T Q_{k+1} H_k = \begin{pmatrix} 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{pmatrix} H_k =: \tilde{H}_k \in \mathbb{R}^{k \times k}. \quad (5.7)$$

Thus, \tilde{H}_k contains just the first k rows of H_k . Since

$$\left(Q_k^T A Q_k\right)^T = Q_k^T A^T Q_k = Q_k^T A Q_k,$$

is a symmetric matrix, \tilde{H}_k is symmetric, too. As in the case of a general matrix, H_k and with that \tilde{H}_k is an upper Hessenberg matrix. From its symmetry, it follows that \tilde{H}_k is even a tridiagonal matrix. Hence, H_k is a tridiagonal matrix, too

$$H_k = \begin{pmatrix} \alpha_1 & \beta_2 & 0 & \cdots & 0 & 0 \\ \beta_2 & \alpha_2 & \beta_3 & \cdots & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \alpha_{k-1} & \beta_k \\ 0 & 0 & 0 & \cdots & \beta_k & \alpha_k \\ 0 & 0 & 0 & \cdots & 0 & \beta_{k+1} \end{pmatrix} \in \mathbb{R}^{(k+1) \times k}.$$

Arnoldi's method simplifies. Using (5.3) and the special form of H_k , one obtains the relation

$$A \underline{q}_k = \beta_k \underline{q}_{k-1} + \alpha_k \underline{q}_k + \beta_{k+1} \underline{q}_{k+1}.$$

From this relation, \underline{q}_{k+1} can be computed. The corresponding algorithm is called Lanczos⁷ algorithm. \square

Algorithm 5.12. Lanczos algorithm – modified Gram–Schmidt variant. Given a symmetric matrix $A \in \mathbb{R}^{n \times n}$ and $\underline{q}_1 \in \mathbb{R}^n$ with $\|\underline{q}_1\|_2 = 1$.

1. $\beta_1 = 0$
2. $\underline{q}_0 = \underline{0}$
3. **for** $j = 1 : m$
4. $\underline{s} = A \underline{q}_j - \beta_j \underline{q}_{j-1}$
5. $\alpha_j = (\underline{s}, \underline{q}_j)$
6. $\underline{s} = \underline{s} - \alpha_j \underline{q}_j$
7. $\beta_{j+1} = \|\underline{s}\|_2$
8. **if** $\beta_{j+1} == 0$

⁷ Cornelius Lanczos (1893 – 1974)

$$\bar{R}_k = \bar{Q}^T H_k = G_k^T G_{k-1}^T \cdots G_2^T G_1^T H_k.$$

Since H_k is tridiagonal, one obtains

$$\bar{R}_k = \begin{pmatrix} r_{11} & r_{12} & r_{13} & 0 & \cdots & 0 \\ 0 & r_{22} & r_{23} & r_{24} & \cdots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & & & & & r_{k-2,k} \\ 0 & & & & & r_{k-1,k} \\ 0 & & & & & r_{k,k} \\ 0 & 0 & 0 & 0 & \cdots & 0 \end{pmatrix} \in \mathbb{R}^{(k+1) \times k}, \quad (5.10)$$

i.e., $r_{ij} = 0$ if $j > i + 2$. A Givens rotation changes only the rows that are involved, i.e., here the two neighboring rows j and $j + 1$. The new rows are linear combinations of the old rows, where just the value zero in the entry $(j + 1, j)$ is produced. A zero value can be converted to a non-zero value only at the entry $(j, j + 2)$, because the entry in $(j + 1, j + 2)$ is generally not zero.

Consider the only interesting case $r^{(k)} \neq 0$, in which the matrix H_k has full rank k . Let R_k be the matrix that consists of the first k rows of \bar{R}_k . The matrix R_k is non-singular since H_k and \bar{Q} have full rank such that \bar{R}_k has rank k . Setting

$$P_k = \begin{pmatrix} p_1 & p_2 & \cdots & p_k \end{pmatrix} := Q_k R_k^{-1} \in \mathbb{R}^{n \times k},$$

one has from $P_k R_k = Q_k$ and due to the special form of R_k , compare (5.10), the recursion

$$\begin{aligned} p_1 &= \frac{q_1}{r_{11}}, \\ p_2 &= \frac{1}{r_{22}} \left(q_2 - r_{12} p_1 \right) \quad \left(\Leftarrow r_{22} p_2 + r_{12} p_1 = q_2 \right), \\ &\vdots \\ p_j &= \frac{1}{r_{jj}} \left(q_j - p_{j-1} r_{j-1,j} - p_{j-2} r_{j-2,j} \right), \quad j = 3, \dots, k. \end{aligned} \quad (5.11)$$

The least squares problem (5.5) can now be rewritten in the form

$$\begin{aligned} \min_{\underline{z} \in \mathbb{R}^k} & \left\| \left\| \underline{r}^{(0)} \right\|_2 \underline{e}_1 - G_1 G_2 \cdots G_k \bar{R}_k \underline{z} \right\|_2^2 \\ &= \min_{\underline{z} \in \mathbb{R}^k} \left\| \left\| \underline{r}^{(0)} \right\|_2 G_k^T \cdots G_2^T G_1^T \underline{e}_1 - \bar{R}_k \underline{z} \right\|_2^2, \end{aligned}$$

because the Euclidean norm is invariant under a multiplication with a unitary matrix. Since the last row of \bar{R}_k vanishes, its Moore⁹-Penrose¹⁰ inverse (pseudo-inverse), see Numerical Mathematics I, is given by

$$\bar{R}_k^+ = \begin{pmatrix} R_k^{-1} & \underline{0} \end{pmatrix} \in \mathbb{R}^{k \times (k+1)}$$

and the solution of the least squares problem is given by

$$\underline{z}^{(k)} = \left\| \underline{r}^{(0)} \right\|_2 \bar{R}_k^+ G_k^T \cdots G_2^T G_1^T \underline{e}_1 = \left\| \underline{r}^{(0)} \right\|_2 R_k^{-1} \begin{pmatrix} G_k^T \cdots G_1^T \underline{e}_1 \end{pmatrix}_{1 \leq i \leq k},$$

⁹ Eliakim Hastings Moore (1862 – 1932)

¹⁰ Roger Penrose, born 1931

where the last index symbolizes that only the first k components of the vectors are taken. Consequently, the iterate with the minimal residual has the form, see (5.6),

$$\begin{aligned}\underline{x}^{(k)} &= \underline{x}^{(0)} + Q_k \underline{z}^{(k)} = \underline{x}^{(0)} + \left\| \underline{r}^{(0)} \right\|_2 Q_k R_k^{-1} \underbrace{\left(G_k^T \cdots G_1^T \underline{e}_1 \right)}_{\in \mathbb{R}^{k+1}} \Big|_{1 \leq i \leq k} \\ &= \underline{x}^{(0)} + \left\| \underline{r}^{(0)} \right\|_2 P_k \left(G_k^T \cdots G_1^T \underline{e}_1 \right) \Big|_{1 \leq i \leq k}.\end{aligned}$$

Since the Givens rotation or reflection G_j^T influences only the components j and $j + 1$ of the vector to which it is applied, the first $(j - 1)$ of its components stay unchanged:

$$\left(G_j^T \cdots G_1^T \underline{e}_1 \right) \Big|_{1 \leq i \leq j-1} = \left(G_{j-1}^T \cdots G_1^T \underline{e}_1 \right) \Big|_{1 \leq i \leq j-1}.$$

It follows that

$$\begin{aligned}\underline{x}^{(k)} &= \underline{x}^{(0)} + \left\| \underline{r}^{(0)} \right\|_2 P_{k-1} \left(G_k^T \cdots G_1^T \underline{e}_1 \right) \Big|_{1 \leq i \leq k-1} + \left\| \underline{r}^{(0)} \right\|_2 \underline{p}_k \left(G_k^T \cdots G_1^T \underline{e}_1 \right) \Big|_{i=k} \\ &= \underbrace{\underline{x}^{(0)} + \left\| \underline{r}^{(0)} \right\|_2 P_{k-1} \left(G_{k-1}^T \cdots G_1^T \underline{e}_1 \right) \Big|_{1 \leq i \leq k-1}}_{=\underline{x}^{(k-1)}} + \left\| \underline{r}^{(0)} \right\|_2 \underline{p}_k \left(G_k^T \cdots G_1^T \underline{e}_1 \right) \Big|_{i=k} \\ &= \underline{x}^{(k-1)} + \left\| \underline{r}^{(0)} \right\|_2 \left(G_k^T \cdots G_1^T \underline{e}_1 \right) \Big|_{i=k} \underline{p}_k.\end{aligned}$$

For computing \underline{p}_k , one needs, see (5.11), \underline{q}_k , \underline{p}_{k-1} , and \underline{p}_{k-2} . The result \underline{p}_k can be stored in place of \underline{p}_{k-2} since \underline{p}_{k-2} is not needed any longer. Hence, the minimization problem can be solved without needing the complete basis of the Krylov subspace.

Together with the short recurrence of the Lanczos algorithm, it is shown that the storage of the basis of $K_k \left(\underline{r}^{(0)}, A \right)$ is not necessary.

The resulting method that computes iterates with minimal residual for symmetric matrices A is called MINRES. MINRES requires to store six arrays: \underline{q}_k , \underline{q}_{k+1} , \underline{s} , \underline{p}_k , \underline{p}_{k-1} , and $\underline{x}^{(k)}$. In contrast to GMRES, the current iterate $\underline{x}^{(k)}$ is known and not only the residual of the current iterate. \square

Remark 5.16. S.p.d. matrices, conjugate residual method. In practice, A is often not only symmetric but also positive definite. In this case, MINRES is seldom used. Even in the context of methods that minimize the residual, there is a more efficient method for s.p.d. matrices called conjugate residual method. \square

Definition 5.17. Conjugate vectors. Let $A \in \mathbb{R}^{n \times n}$ be symmetric and positive definite. The vectors $\underline{x}, \underline{y} \in \mathbb{R}^n$ are called A -orthogonal or (A) -conjugate if

$$\underline{x}^T A \underline{y} = (A \underline{x}, \underline{y}) = 0.$$

If there is no ambiguity, the vectors are called just conjugate. \square

Remark 5.18. Comparison of conjugate residual and conjugate gradient method. The conjugate residual method needs to store only five arrays. It requires in each iteration one matrix-vector product. The memory requirements are one array more than the conjugate gradient method, see Section 6.2. In addition, one has to compute one vector update ($2n$ flops) per iteration more with the conjugate residual method in comparison with the conjugate gradient method. Since both methods need in general a similar number of iterations, the conjugate gradient method is preferred in practice. For this reason, it will be referred to the literature for more details concerning the conjugate residual method. \square